



UNIVERSITY
OF WOLLONGONG
AUSTRALIA

University of Wollongong
Research Online

Faculty of Engineering and Information Sciences -
Papers: Part A

Faculty of Engineering and Information Sciences

2016

Hot enough for you? a spatial exploratory and inferential analysis of North American climate-change projections

Noel A. Cressie

University of Wollongong, ncressie@uow.edu.au

Emily L. Kang

University of Cincinnati

Publication Details

Cressie, N. & Kang, E. L. (2016). Hot enough for you? A spatial exploratory and inferential analysis of North American climate-change projections. *Mathematical Geosciences*, 48 (2), 107-121.

Research Online is the open access institutional repository for the University of Wollongong. For further information contact the UOW Library:
research-pubs@uow.edu.au

Hot enough for you? a spatial exploratory and inferential analysis of North American climate-change projections

Abstract

Climate models have become the primary tools for scientists to project climate-change into the future and to understand its potential impact. Continental-scale General Circulation Models (GCMs) oversimplify the regional climate processes and geophysical features such as topography and land cover. The consequences of local/regional climate change are particularly relevant to natural resource management and environmental-policy decisions, for which Regional Climate Models (RCMs) have been developed. RCMs simulate, for example, three-hourly "weather" over long time periods, from which long-run averages (e.g., over 30 years) are commonly computed to estimate a region's future climate. With anthropogenic forcings incorporated, RCMs provide a means to assess a combination of natural and anthropogenic influences on climate variability. The North American Regional Climate Change Assessment Program ran RCMs into the future, until 2070, for 11,760 contiguous regions, each of which is approximately (Formula presented.) in area. Using the 94,080 temperature changes projected to 2070 for all regions, for two RCMs, and for the four seasons, we present both an exploratory and a Bayesian inferential spatial analysis. Climate-model output is deterministic, but we capture its spatial variability using a hierarchy of conditional probability models. The exploratory Spatial Proportion Over Threshold (SPOT) function and the inferential PRedictive probability Over Threshold (PROT) function are defined and contrasted through videos available online in the Supplementary Materials, showing regions of North America that attain or exceed temperature change thresholds as a function of increasing threshold. The preponderance of our results throughout all regions of North America is one of warming by 2070, usually more (and sometimes much more) than (Formula presented.).

Disciplines

Engineering | Science and Technology Studies

Publication Details

Cressie, N. & Kang, E. L. (2016). Hot enough for you? A spatial exploratory and inferential analysis of North American climate-change projections. *Mathematical Geosciences*, 48 (2), 107-121.

Mathematical Geosciences manuscript No.
(will be inserted by the editor)

Hot Enough for You? A Spatial Exploratory and Inferential Analysis of North American Climate-Change Projections

Noel Cressie · Emily L. Kang

Received: date / Accepted: date

Abstract Climate models have become the primary tools for scientists to project climate change into the future and to understand its potential impact. Continental-scale General Circulation Models (GCMs) oversimplify the regional climate processes and geophysical features such as topography and land cover. Since the consequences of local/regional climate change are particularly relevant to natural-resource management and environmental-policy decisions, Regional Climate Models (RCMs) have been developed. RCMs can simulate three-hourly “weather” over long time periods, from which long-run averages (e.g., over 30 years) are commonly computed to estimate a region’s future climate. With anthropogenic forcings incorporated, RCMs provide a means to assess a combination of natural and anthropogenic influences on climate variability.

N. Cressie
Centre for Environmental Informatics, Building 39c.264,
University of Wollongong, NSW 2522, Australia
Tel.: +61 2 4221 5168
Fax: +61 2 4221 4845
E-mail: ncressie@uow.edu.au

E. Kang
Department of Mathematical Sciences,
University of Cincinnati, Cincinnati OH 45221-0025

The North American Regional Climate Change Assessment Program (NARCCAP) ran regional Climate Models (RCMs) into the future, until 2070, for 11,760 contiguous regions, each of which is approximately $50 \text{ km} \times 50 \text{ km}$ in area. Using the 94,080 temperature changes projected to 2070 for all regions, for two RCMs, and for the four seasons, we present both an exploratory and a Bayesian inferential spatial analysis. Climate-model output is deterministic, but we capture its spatial variability using a hierarchy of conditional probability models. The exploratory SPOT function and the inferential PROT function are defined and contrasted through videos available online in the Supplementary Materials, showing regions of North America that attain or exceed temperature-change thresholds as a function of increasing threshold. The preponderance of our results throughout all regions of North America is one of warming by 2070, usually more (and sometimes much more) than 2°C .

Keywords ESDA · PROT function · spatial hierarchical model · SPOT function · temperature-change projections

1 Introduction

Climate models have become primary tools for scientists to project future climate change and to understand its potential impact. Since the late 1960s, Atmosphere-Ocean General Circulation Models (GCMs) have been developed to simulate the climate over the entire globe. GCMs couple an atmospheric model with an oceanic model to simulate components of the global climate system, such as circulations and forcings. Due to model complexity and limitations of computational resources, GCMs are restricted to generate outputs on coarse spatial scales, typically 200 to 500 km. Additionally, due to their global perspective, GCMs oversimplify the regional climate processes and geophysical features such as topography and land cover. Since local/regional climate

1 effects are more relevant to natural-resource management and environmental-
2 policy decisions, Regional Climate Models (RCMs) have been developed to
3 produce high-resolution outputs on scales of 50 km and smaller. Nevertheless,
4 RCMs need initial conditions and time-dependent boundary conditions, which
5 are typically provided by a GCM; this is sometimes referred to as “dynamic
6 downscaling” of the GCM outputs (e.g., Fennessy and Shukla 2000; Xue et al.
7 2007).

8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
Essentially, GCMs and RCMs are a series of discretized differential equa-
tions that attempt to represent physical relationships such as the flows of
energy and water within and between the atmosphere, oceans, land, sea ice,
and so forth. Using differential equations that describe the physical dynamics,
RCMs can simulate three-hourly “weather” over long time periods and gener-
ate a vast array of outputs, from which the long-run average is commonly used
as a summary of how a climate model approximates a region’s climate. With
anthropogenic forcings incorporated, climate models can be run under different
scenarios (e.g., various CO₂ levels). Consequently, natural and anthropogenic
influences on climate variability can be assessed.

31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
GCMs and RCMs are complicated to build, and they are deterministic (i.e.,
two runs of the same model under identical conditions produce identical out-
puts). Nevertheless, there are various sources of uncertainty that accompany
such models. For example, there may be uncertainty in assumptions about
interaction between atmospheric circulation and orography, about discretiza-
tion, or about parameterizations of the physical-forcing processes. To obtain
a better understanding of such uncertainties, climate scientists carry out ex-
periments with multiple runs of multiple models. In this article, we consider
climate-model output from the North American Regional Climate Change As-
sessment Program (NARCCAP; Mearns et al. 2009), and we concentrate on
RCMs with the same initial and boundary conditions (supplied by the same

GCM) and the same greenhouse-gas (GHG) forcing (supplied by a scenario that matches most closely current GHG emissions; Nakicenovic et al. 2000).

NARCCAP is an international program whose goal is to produce high-resolution climate simulations for current and future time periods. Thus, there is an opportunity to investigate spatial variability of regional-scale projections of future climate and to generate temperature-change scenarios for use in impacts research. NARCCAP produces this high-resolution (approximately 50 km) climate-output data for six RCMs built for the North American region that includes the US, Canada, northern Mexico, and the oceans nearby (Mearns et al. 2009).

NARCCAP Phase I explores the variability in RCM outputs for the current period, where the six RCMs were run with common boundary conditions provided by the NCEP-DOE Reanalysis II data (e.g., Kanamitsu et al. 2002). NARCCAP Phase II involves not only multiple RCMs but also runs with different boundary conditions provided by different GCMs.

In Phase II, these six RCMs are coupled with four different GCMs and were run not only for the current period (1971-2000) but also for a future period (2041-2070). Thus, temperature-change projections are available from the Phase II experiment. In both Phase I and Phase II, the same greenhouse-gas emissions scenario (SRES A2; Nakicenovic et al. 2000) were used.

A number of different aspects of the NARCCAP data have already been analyzed (e.g., Kaufman and Sain 2010; Salazar et al. 2011; Kang et al. 2012; Kang and Cressie 2013). In this article, attention is paid to the seasonal climate projections, particularly the Boreal winter, reasoning that projected warming over North America might be most clearly seen in the months of December, January, and February; see Kang and Cressie (2013).

In this article, we consider a subset of the output from the Phase II runs, namely the average surface yearly temperatures and the average sur-

face Boreal-winter (December, January, and February) temperatures, for the current period (1971-2000) and for the future period (2041-2070). We analyze the output produced by two RCMs, namely CRCM and RCM3, with the same GCM, namely CGCM3, providing the boundary conditions; for further details on these and other climate models used in NARCCAP, see Kang and Cressie (2013). The RCM outputs for the current and future periods were given on an approximately $50 \text{ km} \times 50 \text{ km}$ NARCCAP grid of 98×120 pixels, resulting in 11,760 NARCCAP pixels. These are “big data”; there are a total of 5.6 million surface temperature values used in the exploratory and inferential analyses presented in this article.

Section 2 gives exploratory spatial data analyses (ESDA) of the subset of Phase II output under consideration. The SPOT function is introduced and illustrated on the yearly and winter temperature-change data. Section 3 presents a Bayesian hierarchical spatial model of the NARCCAP data used in Sect. 2; and inferential analysis is given in Sect. 4, where the PROT function is introduced. A discussion and various conclusions are given in Sect. 5. Videos showing the exploratory SPOT function and the inferential PROT function are provided online in the Supplementary Materials.

2 Exploratory Spatial Data Analysis and the SPOT Function

Consider the current and future 30-year-averaged seasonal surface temperature fields from the i -th RCM for the j -th Boreal season, where $i = 1, 2$, and $j = 1$ (spring: March, April, May), $j = 2$ (summer: June, July, August), $j = 3$ (autumn: September, October, November), $j = 4$ (winter: December, January,

February). Define

$$Z_{ij}^{\text{current}}(\mathbf{s}_l) \equiv (1/30) \sum_{t=1971}^{2000} Z_{ij}(\mathbf{s}_l; t)$$

$$Z_{ij}^{\text{future}}(\mathbf{s}_l) \equiv (1/30) \sum_{t=2041}^{2070} Z_{ij}(\mathbf{s}_l; t),$$

where $l = 1, \dots, 11760$ represents the index of the l -th pixel on the 98×120 NARCCAP grid that is superimposed on North America. Here, $Z_{ij}(\mathbf{s}_l; t)$ is the corresponding region's average surface temperature for the i -th RCM and the j -th season in year t . In all, there are 5.6 million spatio-temporal Z -values used in the two definitions given above.

In this article, we are particularly interested in temperature change for the winter ($j = 4$) and for the whole year (averaged over $j = 1, 2, 3, 4$). Further, we mostly consider the output of the RCMs averaged over the models $i = 1, 2$, although some discussion of the between-model variability is given in Sect. 5; it will be seen to be small (about 5% of the total variability). Define the temperature-change projections,

$$D^{wi}(\mathbf{s}_l) \equiv (1/2) \sum_{i=1}^2 Z_{i4}^{\text{future}}(\mathbf{s}_l) - (1/2) \sum_{i=1}^2 Z_{i4}^{\text{current}}(\mathbf{s}_l) \quad (1)$$

$$D^{yr}(\mathbf{s}_l) \equiv (1/8) \sum_{i=1}^2 \sum_{j=1}^4 Z_{ij}^{\text{future}}(\mathbf{s}_l) - (1/8) \sum_{i=1}^2 \sum_{j=1}^4 Z_{ij}^{\text{current}}(\mathbf{s}_l). \quad (2)$$

These represent projections of (30-year-averaged) North American surface temperature changes, projected out to 2070. Each of (1) and (2) involves 11,760 data defined on the NARCCAP grid. This spatial context means that visualizations through mapping can (and should) be part of any exploratory data analysis of temperature-change projections obtained from NARCCAP.

It can be seen from Fig. 1 that the temperature changes are consistently positive for both winter output (upper-left panel) and yearly output (upper-

right panel). That is, it is projected to be considerably warmer by 2070 over the entire North American region. Generally speaking, the warming effect is stronger over land compared to that over ocean. It is also apparent that the warming effect during the winter in the northern part of the North America region is particularly strong, especially in the Hudson Bay area.

These visualizations represent “map views,” which can be enhanced by a “multivariate view,” namely an x - y plot of winter output versus yearly output. The lower-left panel of Fig. 1 shows a plot of the 11760 points $\{(D^{yr}(\mathbf{s}_l), D^{wi}(\mathbf{s}_l)) : l = 1, \dots, 11760\}$. As expected, the winter temperature changes are more extreme than the yearly ones, since the yearly changes are averages over all seasons. Obviously, people and communities live through seasons, and hence it is highly relevant to analyze seasonal behavior as well as yearly behavior.

Figure 1 here

To quantify these impressions, a function is introduced that plots the spatial proportion (of pixels) over a threshold as a function of that threshold. This SPOT (Spatial Proportion Over Threshold) function is defined as follows.

$$T^{wi}(k) \equiv (1/n) \sum_{l=1}^n I(D^{wi}(\mathbf{s}_l) > k); \quad -\infty < k < \infty \quad (3)$$

$$T^{yr}(k) \equiv (1/n) \sum_{l=1}^n I(D^{yr}(\mathbf{s}_l) > k); \quad -\infty < k < \infty, \quad (4)$$

where $n = 11760$ is the number of spatial regions being considered. A very powerful technique in exploratory spatial data analysis is to “link and brush,” namely to “paint” on a map of North America all those regions $\{\mathbf{s}_l\}$ whose projected temperature change exceeds the threshold k . Fig. 2 shows a video of the “SPOT view” and “map view,” linked, as k increases from 1.0°C to 7.6°C in units of 0.2°C (video available online in the Supplementary Materials). Both

$T^{wi}(\cdot)$ and $T^{yr}(\cdot)$ are shown. The linking of the SPOT view and the map view shows where and what proportion of pixels exceed a range of thresholds.

Figure 2 here

The role of exploratory spatial data analysis is to generate ideas about specific sources of variability, which can be followed up by fitting statistical models and carrying out statistical inference. Here, spatial variation is an obvious source of variability, which we have already seen in Figs. 1 and 2. In the next section, we quantify the variation of the RCMs' outputs with a stochastic model and use it to address the questions, "Is climate change real?" and "What are the projected temperature changes along with measures of their certainty?" Critically, this is done with a paradigm that expresses uncertainties through probability distributions.

3 Spatial Hierarchical Model of Temperature Change by 2070

Recall the definitions (1) and (2) obtained from the output of a collection of (here two) RCMs. It is not expected that in 2070, temperatures will increase exactly by any of these amounts. While the output of an RCM is deterministic, there is uncertainty about how continuous-time, continuous-space partial differential equations should be discretized onto grids, along with where and how forcing terms (e.g., greenhouse gases) are introduced. Thus, an RCM's output could be thought of as a single realization of many possible outputs, which are governed by a probability distribution (e.g., Kaufman and Sain 2010; Sain et al. 2011; Salazar et al. 2011; Kang et al. 2012; Kang and Cressie 2013). It should be noted that geostatistics takes exactly the same approach when making inferences on an ore body, even though the mineralization of a given ore body was a unique event.

3.1 Introduction to Hierarchical Statistical Modeling

A hierarchical statistical model is one where the model can be broken down into at least two levels: The data model and the process model. A Bayesian hierarchical statistical model involves three levels: The data model, the process model and, additionally, the parameter (or prior) model. When these models are multiplied together, they form the joint distribution of the data and the process/parameters (e.g., Berliner 1996). The data model describes the conditional probability distribution of the data, given parameters and an unobserved (hidden) process. The process model describes the conditional probability distribution of the hidden process given its parameters. The parameter model puts a “prior” distribution on the parameters themselves.

In the application presented here, the data model describes the long-run average differences between future and current climate-model runs, where the hidden climate process is made up of the projected temperature changes by season and RCM. The process model incorporates the Spatial Random Effects (SRE) model (Cressie and Johannesson 2008), which is an effective way to reduce the dimensionality of the problem from approximately 100,000 to less than 100. Prior distributions are assigned to parameters, which constitutes the parameter model. The ultimate goal is to obtain the posterior distribution, namely the joint distribution of the unknowns in the hierarchical statistical model (i.e., process and parameters), given the data. The predictive distribution is simply the marginal distribution of the unknown process given the data. Using Bayes’ Theorem, the posterior distribution is proportional to the product of the data, process, and parameter models. Simulation procedures, such as Markov chain Monte Carlo (MCMC) methods, are used here to obtain the predictive distribution of any part of the process (given the data).

Further details on the data/process and data/process/parameters hierarchical framework can be found in Cressie and Wikle (2011, Chap. 2).

There are several advantages to using a hierarchical statistical approach. First, non-hierarchical models with just a few parameters generally do not fit the data well. Moreover, non-hierarchical models with many parameters may fit the data well but tend to “over-fit” and may not be useful for predictive purposes. In contrast, hierarchical statistical models can often fit the data well with just a few parameters, and they also do well for predicting the hidden process (at both observed and unobserved parts of the process).

3.2 Data Model

This part of the hierarchical statistical model usually incorporates the component of variability due to measurement error. But, it can also capture other sources of variability extraneous to the hidden process of interest, such as spatio-temporal interactions. Recall the definition of $Z_{ij}(\mathbf{s}_l; t)$ as the l -th region’s average surface temperature for the i -th RCM and the j -th season in year t .

For $t = 2041, \dots, 2070$, write

$$Z_{ij}(\mathbf{s}_l; t) - Z_{ij}(\mathbf{s}_l; t - 70) = Y_{ij}(\mathbf{s}_l) + e_{ij}(\mathbf{s}; t), \quad (5)$$

where $Y_{ij}(\cdot)$ is the purely spatial process of temperature change by 2070, and $e_{ij}(\cdot; t)$ is an independent smaller-scale component of variation capturing spatio-temporal interaction. It is assumed in (5) that Y_{ij} carries all the spatial dependence in the temperature-change field and that e_{ij} represents small errors where the spatial and temporal dependencies are negligible. Consequently,

we assume that e_{ij} has mean zero and, independently for $t = 2041, \dots, 2071$,

$$e_{ij}(\cdot; t) \sim \text{Gau}(0, V_{ij}(\cdot)\sigma_e^2).$$

Define the 30-year average temperature difference,

$$D_{ij}(\mathbf{s}_l) \equiv (1/30) \sum_{t=2041}^{2070} \{Z_{ij}(\mathbf{s}_l; t) - Z_{ij}(\mathbf{s}_l; t - 70)\}, \quad (6)$$

and hence, from (5), $D_{ij}(\cdot)$ can be written as

$$D_{ij}(\cdot) = Y_{ij}(\cdot) + \varepsilon_{ij}(\cdot), \quad (7)$$

where the error term ε_{ij} has no spatial dependence, and

$$\varepsilon_{ij}(\mathbf{s}_l) \sim \text{Gau}(0, V_{ij}(\mathbf{s}_l)\sigma_\varepsilon^2),$$

for $\sigma_\varepsilon^2 = (1/30)\sigma_e^2$. The spatial heterogeneity of the error term is captured by $V_{ij}(\cdot)$, which is given by

$$V_{ij}(\mathbf{s}_l) \equiv (1/29) \sum_{t=2041}^{2070} \{Z_{ij}(\mathbf{s}_l; t) - Z_{ij}(\mathbf{s}_l; t - 70) - D_{ij}(\mathbf{s}_l)\}^2.$$

Equation (7) is the data model. The goal is to filter out the error term ε_{ij} and to make inference on the hidden spatial process Y_{ij} ; this is the formulation developed by Kang and Cressie (2013).

In this article, we are interested in temperature changes for the winter and for the whole year. From (1), (2), and (7),

$$D^{wi}(\mathbf{s}_l) = Y^{wi}(\mathbf{s}_l) + (1/2) \left(\sum_{i=1}^2 \varepsilon_{i4}(\mathbf{s}_l) \right) \quad (8)$$

$$D^{yr}(\mathbf{s}_l) = Y^{yr}(\mathbf{s}_l) + (1/8) \left(\sum_{i=1}^2 \sum_{j=1}^4 \varepsilon_{ij}(\mathbf{s}_l) \right), \quad (9)$$

for $l = 1, \dots, 11760$; and in (8) and (9),

$$Y^{wi}(\cdot) \equiv (1/2) \sum_{i=1}^2 Y_{i4}(\cdot); \quad Y^{yr}(\cdot) \equiv (1/8) \sum_{i=1}^2 \sum_{j=1}^4 Y_{ij}(\cdot). \quad (10)$$

Inference will be on the two hidden spatial processes, Y^{wi} and Y^{yr} , which from (10) are defined in terms of Y_{ij} given in (7).

3.3 Process and Parameter Models

The variability in the temperature-change process Y_{ij} is due to climate-model differences, seasonal differences, and spatial variability. This can be described through the decomposition

$$Y_{ij}(\cdot) = \mu(\cdot) + a_i(\cdot) + b_j(\cdot) + (ab)_{ij}(\cdot), \quad (11)$$

where $\mu(\cdot)$ represents the baseline temperature change from which RCM and seasonal variation will be assessed; $a_i(\cdot)$ is the extra component due to the i -th RCM; $b_j(\cdot)$ is the extra component due to the j -th season; and the terms $\{(ab)_{ij}(\cdot)\}$ capture the RCM-season interaction and typically exhibit smaller-scale variability. Notice that (11) assumes there is no fine-scale variability (i.e., no nugget effect) since the 30-year averaging of climate-model output results in a spatially smooth process.

All the components on the right-hand side of (11) are spatial processes, each with $n = 11760$ elements. The spatial covariance matrices of these processes are 11760×11760 and generally too large to be inverted. Kang and Cressie (2013) proposed the following dimension reduction,

$$Y_{ij}(\cdot) = \mu(\cdot) + \mathbf{S}(\cdot)' \boldsymbol{\eta}_{ij}, \quad (12)$$

where $\mu(\cdot)$ is a deterministic trend that may depend on covariates through $\mu(\cdot) = \mathbf{x}(\cdot)'\boldsymbol{\beta}$; $\mathbf{S}(\cdot)$ is a known r -dimensional vector defined by r spatial basis functions such that $r \ll n$; and the r -dimensional, zero-mean random vector of coefficients, $\boldsymbol{\eta}_{ij}$, is decomposed into independent zero-mean components according to

$$\boldsymbol{\eta}_{ij} = \boldsymbol{\delta}_i + \boldsymbol{\gamma}_j + \boldsymbol{\zeta}_{ij}. \quad (13)$$

For $i = 1, 2$, Kang and Cressie (2013) fitted the model defined by the random effects due to

RCM: $\boldsymbol{\delta}_i \sim \text{Gau}(\mathbf{0}, \mathbf{K}_1)$

Season: $\boldsymbol{\gamma}_j \sim \text{Gau}(\mathbf{0}, \mathbf{K}_2), j = 1, 3; \quad \boldsymbol{\gamma}_j \sim \text{Gau}(\mathbf{0}, \mathbf{K}_3), j = 2, 4$

Interaction: $\boldsymbol{\zeta}_{ij} \sim \text{Gau}(\mathbf{0}, \mathbf{K}_4), j = 1, 3; \quad \boldsymbol{\zeta}_{ij} \sim \text{Gau}(\mathbf{0}, \mathbf{K}_5), j = 2, 4,$

which allows the covariance matrices to differ between the mild spring and autumn seasons and the more extreme summer and winter seasons. This modeling decision was based on observing that the raw differences in summer and winter, respectively $\{D_{i2}(\cdot)\}$ and $\{D_{i4}(\cdot)\}$, exhibited roughly similar extreme behavior, albeit in different regions. Thus we use the same covariance matrices for the summer and winter seasons. We assume that the baseline temperature change $\mu(\cdot)$ is constant (i.e., we use one covariate, equal to 1 everywhere). The spatial basis functions are those chosen by Kang and Cressie (2013), which capture multi-resolution scales of variability, as well as elevation and presence/absence of pixels on land, the Great Lakes, Hudson Bay, and coastlines. These latter spatial basis functions pay special attention to regions with a lot of energy exchange between land and water. The final tally is $r = 85$ spatial basis functions.

Equations (12) and (13) represent the process model, and (12) is called a Spatial Random Effects (SRE) model; see Cressie and Johannesson (2008). Finally, the parameter models that include prior distributions for $\mathbf{K}_1, \dots, \mathbf{K}_5$, are set out in full detail in Kang and Cressie (2013, Sect. 3.3), to which we refer the interested reader.

3.4 Predictive Distributions of 11,760 Temperature-Change Values

Very simply, the data model is the distribution of D given Y (first level), the process model is the distribution of Y (second level), and the predictive distribution is the distribution of

$$[Y \text{ given } D] \propto [D \text{ given } Y] \times [Y],$$

by Bayes' Theorem. The extra level in a Bayesian hierarchical model is the parameter model, for which the generalization of Bayes' Theorem is straightforward. This third level does not change the need for obtaining a predictive distribution of Y given D .

Fundamentally, inference on the process is conditional on the data. In the NARCCAP context, Y is made up of 94,080 process values $\{Y_{ij}(\cdot)\}$, D is made up of 94,080 data values $\{D_{ij}(\cdot)\}$, and Bayes' Theorem is used to obtain the predictive distribution of $\{Y_{ij}(\cdot)\}$ given $\{D_{ij}(\cdot)\}$. However, with the large number of data values and process values, it is generally computationally infeasible to obtain the predictive distribution, and so some form of dimension reduction is needed.

Randomness appears in the data model (7) very simply through independent error, and it appears in the process model (12) and (13) through the r -dimensional ($r = 85$) vectors $\{\boldsymbol{\eta}_{ij} : i = 1, 2; j = 1, 2, 3, 4\}$. Critically, if the $85 \times 2 \times 4 = 680$ process values in (13) can be inferred, then all of

$\{Y_{ij}(\cdot) : i = 1, 2; j = 1, 2, 3, 4\}$ can be inferred. Consequently, dimension-reduced inference is obtained through the predictive distribution of

$$\{\boldsymbol{\eta}_{ij}\} \text{ given } \{D_{ij}(\mathbf{s}_l) : i = 1, 2; j = 1, 2, 3, 4; l = 1, \dots, 11760\}.$$

Kang and Cressie (2013) give the details of an MCMC algorithm that is used to sample from the predictive distribution of $\{Y_{ij}(\cdot)\}$ given $\{D_{ij}(\cdot)\}$ based on the dimension reduction afforded by the SRE model in (12). Briefly, we ran two parallel Markov chain Monte Carlo (MCMC) chains for 12,500 iterations each, and we discarded the first 2,500 to allow for burn-in. The result was a sample of size 20,000 from the predictive distribution of

$$Y^{wi}(\cdot) \text{ given } \{D_{ij}(\cdot)\}; \quad Y^{yr}(\cdot) \text{ given } \{D_{ij}(\cdot)\}.$$

All the computations were carried out in Matlab on a dual core 2.88 GHz Intel Xeon processor with 96 Gb of memory running Linux; in total, they took 40 CPU hours. All inferences presented in this article are obtained from such predictive distributions. In the next section, inference on an individual region's temperature change as well as inferential analogues to the exploratory SPOT functions (Sect. 2) are presented.

4 Spatial Inference and the PROT Function

In Sect. 3.4, it was explained how MCMC samples from the predictive distribution of $\{Y_{ij}(\cdot)\}$ given the data $\{D_{ij}(\cdot)\}$ can be used to obtain samples from the predictive distribution of $Y^{wi}(\cdot)$ and $Y^{yr}(\cdot)$ given the data $\{D_{ij}(\cdot)\}$. For example, consider the pixel \mathbf{s}_0 containing the National Center for Atmospheric Research (NCAR) Mesa Lab in Boulder, CO, where the NARCCAP project was conducted. The “NCAR” pixel's predictive distributions, of $Y^{wi}(\mathbf{s}_0)$ given

the data and of $Y^{yr}(\mathbf{s}_0)$ given the data, are shown in Fig. 3. Notice that the predictive distribution of temperature change is centered at about 2°C for the winter season and at about 2.7°C for the whole year. In both situations, there is a predictive probability of 1.000 that climate change will be greater than 1.5°C by 2070.

Figure 3 here

The predictive distribution of any function of $\{Y_{ij}(\cdot)\}$, not just a linear function, can be obtained immediately from the MCMC samples. For the l -th pixel located at \mathbf{s}_l , define the functions

$$I(Y^{wi}(\mathbf{s}_l) > k); \quad I(Y^{yr}(\mathbf{s}_l) > k), \quad (14)$$

for $-\infty < k < \infty$. Notice that the indicator function $I(\cdot)$ in (14) is nonlinear and has a range of $\{0, 1\}$. This function is of interest because beyond certain temperature thresholds, insects may thrive, crops may fail, and native plant species may relocate.

One convenient summary of the predictive distribution is its mean. In terms of the data D and the process Y introduced in Sect. 3.4, suppose one wishes to make inference on the (possibly nonlinear) function $g(Y)$. Then an often-used predictor is the predictive mean

$$\hat{g} \equiv E(g(Y)|D),$$

which is obviously a function of the data. In the case of $g(\cdot)$ defined by (14), the predictors for the winter season and for the whole year are, in obvious

notation,

$$\hat{g}^{wi}(\mathbf{s}_l; k) \equiv \Pr(Y^{wi}(\mathbf{s}_l) > k | \{D_{ij}(\cdot)\}) \quad (15)$$

$$\hat{g}^{yr}(\mathbf{s}_l; k) \equiv \Pr(Y^{yr}(\mathbf{s}_l) > k | \{D_{ij}(\cdot)\}), \quad (16)$$

for $l = 1, \dots, 11760$. Definitions (15) and (16) are the predictive probabilities over a threshold as a function of the threshold. This PROT (PRedictive probability Over Threshold) function is formally defined as follows: For $l = 1, \dots, 11760$,

$$P^{wi}(k; \mathbf{s}_l) \equiv \Pr(Y^{wi}(\mathbf{s}_l) > k | \{D_{ij}(\cdot)\}); \quad -\infty < k < \infty \quad (17)$$

$$P^{yr}(k; \mathbf{s}_l) \equiv \Pr(Y^{yr}(\mathbf{s}_l) > k | \{D_{ij}(\cdot)\}); \quad -\infty < k < \infty. \quad (18)$$

Notice that the PROT functions depend on the location \mathbf{s}_l in the spatial domain, unlike the SPOT functions. Hence, for each threshold k , a choropleth map can be made of the predictive probabilities that projected temperature changes will exceed k . This distinction between the two types of plots is huge, and it emphasizes the incredible benefit of spatial modeling and its associated spatial inference. Fig. 4 shows a video of maps based on the PROT functions, $\{P^{wi}(k; \mathbf{s}_l) : l = 1, \dots, 11760\}$ and $\{P^{yr}(k; \mathbf{s}_l) : l = 1, \dots, 11760\}$, as k increases from 1.0°C to 6.8°C in increments of 0.2°C (video available online in the Supplementary Materials).

Figure 4 here

The PROT values in (17) and (18) are generated for each NARCCAP pixel, but they could also be generated for a coarser resolution of the NARCCAP region. Now inference would be on average temperature change over, say, a coarser-resolution block B. For example, for winter, the (nonlinear) function

of interest is

$$I \left(\left(\sum_{\mathbf{s} \in B} Y^{wi}(\mathbf{s}) / \sum_{\mathbf{s} \in B} 1 \right) > k \right).$$

Then inference proceeds exactly as before, based on MCMC samples from the predictive distribution of $\{\boldsymbol{\eta}_{ij}\}$ given the data.

5 Discussion and Conclusions

This article has given complementary analyses that show North America's climate will be considerably warmer by 2070. A climate-change skeptic might argue that the two RCMs give different results, and so neither can be trusted; we use statistical summaries of variability to quantify the consistency of the two climate models. Fig. 5 shows two maps of the differences of the two model outputs, one for winter and one for the whole year. That is,

$$D_{14}(\cdot) - D_{24}(\cdot); \quad (1/4) \sum_{j=1}^4 D_{1j}(\cdot) - (1/4) \sum_{j=1}^4 D_{2j}(\cdot)$$

are maps for winter and the whole year, respectively. The scale on the maps, in degrees C, indicates very little difference between the two RCMs.

Figure 5 here

To formalize this impression, consider the following measure of between-model variability, defined in terms of the hidden process $\{Y_{ij}(\cdot)\}$.

$$R^{yr} \equiv \frac{\left\{ (1/2) \sum_{\mathbf{s} \in D} (Y_1^{yr}(\mathbf{s}) - Y_2^{yr}(\mathbf{s}))^2 / \sum_{\mathbf{s} \in D} 1 \right\}^{1/2}}{\left| \sum_{\mathbf{s} \in D} Y^{yr}(\mathbf{s}) / \sum_{\mathbf{s} \in D} 1 \right|}, \quad (19)$$

where $Y_i^{yr}(\cdot) \equiv (1/4) \sum_{j=1}^4 Y_{ij}(\cdot)$ for the i -th RCM ($i = 1, 2$). The expression (19) is a nonlinear function of the hidden process and expresses the between-model variability as a fraction of the average temperature change over the

whole spatial domain; it is reminiscent of the coefficient of variation in non-spatial contexts.

Since R^{yr} is a function of $\{Y_{ij}(\cdot)\}$, one can obtain the predictive distribution of R^{yr} , given the data, by MCMC (Sect. 4). In this case, the predictive mean of R^{yr} is 4.55% with a two-sided 95% prediction interval of (4.43%, 4.67%). Analogous to (19), R^{wi} is a measure of the between-model variability for winter, and a predictive analysis shows that it is likewise small.

Apart from our results about future climate, which put high probability on much of North America exceeding 2°C by 2070, this article has demonstrated how formal (Bayesian) spatial statistical inference can be carried out on deterministic climate-model outputs. The inferential PROT function shows the important effect of spatial-variability modeling and shrinkage, when compared to the exploratory SPOT function, as demonstrated in videos available online in the Supplementary Materials.

Climate changes produced by RCMs could be used in agriculture to forecast crop-variety yields, in ecology to forecast bird-migration patterns, and in emergency services to forecast bushfire danger. A sobering caveat to these statistical analyses is that there may be common errors made across all RCMs but these errors are unknown; that is, there are “unknown unknowns.” The result would be outputs, $\{Z_{ij}^{\text{future}}(\cdot)\}$, that give a biased view of what future climate will be like; current analyses, statistical and otherwise, account for the known unknowns.

Acknowledgements This research was partially supported by the NASA’s Earth Science Technology Office through its Advanced Information Systems Technology Program. We wish to thank the North American Regional Climate Change Assessment Program (NARCCAP) for providing the data used in this article. NARCCAP is funded by NSF, DoE, NOAA, and EPA’s Office of Research and Development. Many thanks go to Andrew Holder for his

assistance in preparation of this article, to the referees for their excellent suggestions, and to the Editor for his vision regarding new directions for mathematical geosciences.

References

- Berliner LM (1996) Hierarchical Bayesian time series models. In: Hanson K, Silver R (eds) Maximum entropy and Bayesian methods, Kluwer Academic Publishers, Dordrecht, NL, pp 15–22
- Cressie N, Johannesson G (2008) Fixed rank kriging for very large spatial data sets. *Journal of the Royal Statistical Society, Series B* 70:209–226
- Cressie N, Wikle CK (2011) *Statistics for Spatio-Temporal Data*. John Wiley and Sons, Hoboken, NJ
- Fennessy MJ, Shukla J (2000) Seasonal prediction over North America with a regional model nested in a global model. *Journal of Climate* 13:2605–2627
- Kanamitsu M, Ebisuzaki W, Woollen J, Yang SK, Hnilo JJ, Fiorino M, Potter GL (2002) NCEP-DOE AMIP-II Reanalysis (R-2). *Bulletin of the American Meteorological Society* 83:1631–1644
- Kang EL, Cressie N (2013) Bayesian hierarchical ANOVA of regional climate-change projections from NARCCAP Phase II. *International Journal of Applied Earth Observation and Geoinformation* 22:3–15
- Kang EL, Cressie N, Sain SR (2012) Combining outputs from the NARCCAP regional climate models using a Bayesian hierarchical model. *Journal of the Royal Statistical Society, Series C (Applied Statistics)* 61:291–313
- Kaufman CG, Sain SR (2010) Bayesian ANOVA modeling using Gaussian process prior distributions. *Bayesian Analysis* 5:123–150
- Mearns LO, Gutowski W, Jones R, Leung R, McGinnis S, Nunes A, Qian Y (2009) A regional climate change assessment program for North America. *Eos, Transactions American Geophysical Union* 90(36):311

1 Nakicenovic N, Alcamo J, Davis G, de Vries B, Fenhann J, Gaffin S, Gregory
2 K, Grubler A, Jung TY, Kram T, La Rovere EL, Michaelis L, Mori S,
3 Morita T, Pepper W, Pitcher HM, Price L, Riahi K, Roehrl A, Rogner HH,
4 Sankovski A, Schlesinger M, Shukla P, Smith SJ, Swart R, van Rooijen S,
5 Victor N, Dadi Z (2000) Special Report on Emissions Scenarios: a special
6 report of Working Group III of the Intergovernmental Panel on Climate
7 Change. Tech. rep., Environmental Molecular Sciences Laboratory, Pacific
8 Northwest National Laboratory, Richland, WA, USA
9

10 Sain SR, Furrer R, Cressie N (2011) A spatial analysis of multivariate output
11 from regional climate models. *Annals of Applied Statistics* 5:150–175
12

13 Salazar ES, Finley A, Hammerling D, Steinsland I, Wang X, Delamater P
14 (2011) Comparing and blending regional climate model predictions for the
15 American southwest. *Journal of Agricultural, Biological, and Environmental*
16 *Statistics* 16:586–605
17

18 Xue YK, Vasic R, Janjic Z, Mesinger F, Mitchell KE (2007) Assessment
19 of dynamic downscaling of the continental US regional climate using the
20 Eta/SSiB regional climate model. *Journal of Climate* 20:4172–4193
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Figure Captions

Figure 1: ESDA of temperature change in degrees C in North America. Upper left: Map of $\{D^{wi}(\mathbf{s}_l) : l = 1, \dots, 11760\}$, where the top of the scale refers to temperature changes of 5°C and above. Upper right: Map of $\{D^{yr}(\mathbf{s}_l) : l = 1, \dots, 11760\}$. Lower left: Plot of $\{(D^{yr}(\mathbf{s}_l), D^{wi}(\mathbf{s}_l)) : l = 1, \dots, 11760\}$, where the units on both axes are degrees Celsius.

Figure 2: Video showing linked views of temperature change in degrees C in North America. Upper left and upper right: SPOT function $T^{wi}(k)$ as a function of k linked to the map of $\{I(D^{wi}(\mathbf{s}_l) > k) : l = 1, \dots, 11760\}$. Lower left and lower right: SPOT function $T^{yr}(k)$ as a function of k linked to the map of $\{I(D^{yr}(\mathbf{s}_l) > k) : l = 1, \dots, 11760\}$. Units on the horizontal axis of both SPOT functions are degrees Celsius. The video is created by varying k from 1.0°C to 7.6°C in steps of 0.2°C .

Figure 3: The NARCCAP pixel \mathbf{s}_0 that contains NCAR's Mesa Lab is featured. Upper panel: Predictive distribution of $Y^{wi}(\mathbf{s}_0)$ given the data $\{D_{ij}(\cdot)\}$. Lower panel: Predictive distribution of $Y^{yr}(\mathbf{s}_0)$ given the data $\{D_{ij}(\cdot)\}$. For both plots, the vertical axis shows counts (out of 20,000) and units on the horizontal axis are degrees Celsius.

Figure 4: Video showing inferential analysis of temperature change in degrees C in North America. Left panel: Map of PROT function $\{P^{wi}(k; \mathbf{s}_l) : l = 1, \dots, 11760\}$. Right panel: Map of PROT function $\{P^{yr}(k; \mathbf{s}_l) : l = 1, \dots, 11760\}$. Units on the horizontal axis of both PROT functions are degrees Celsius. The video is created by varying k from 1.0°C to 6.8°C in steps of 0.2°C .

Figure 5: Between-model variability of temperature change in degrees C in North America. Left panel: Map of $\{D_{14}(\mathbf{s}_l) - D_{24}(\mathbf{s}_l) : l = 1, \dots, 11760\}$ for the winter season. Right panel: Map of $\{(1/4) \sum_{j=1}^4 D_{1j}(\mathbf{s}_l) - (1/4) \sum_{j=1}^4 D_{2j}(\mathbf{s}_l) : l = 1, \dots, 11760\}$ for the whole year. The scale for both maps are in units of degrees Celsius.

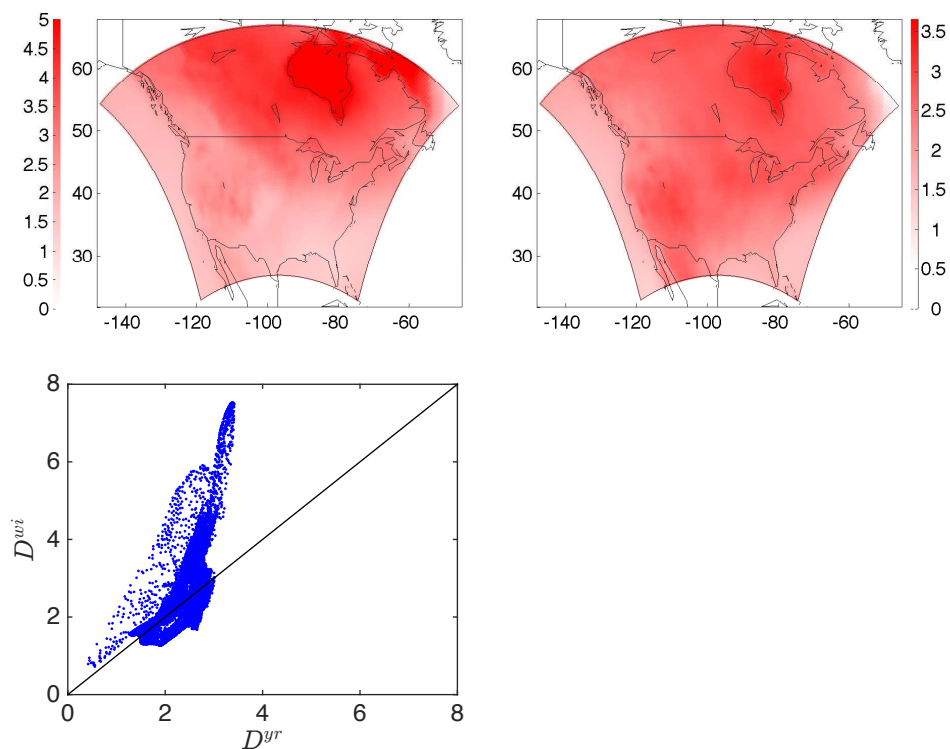


Fig. 1: ESDA of temperature change in degrees C in North America. Upper left: Map of $\{D^{wi}(\mathbf{s}_l) : l = 1, \dots, 11760\}$, where the top of the scale refers to temperature changes of 5°C and above. Upper right: Map of $\{D^{yr}(\mathbf{s}_l) : l = 1, \dots, 11760\}$. Lower left: Plot of $\{(D^{yr}(\mathbf{s}_l), D^{wi}(\mathbf{s}_l)) : l = 1, \dots, 11760\}$, where the units on both axes are degrees Celsius.

(Figure 2 video)

Fig. 2: Video showing linked views of temperature change in degrees C in North America. Upper left and upper right: SPOT function $T^{wi}(k)$ as a function of k linked to the map of $\{I(D^{wi}(\mathbf{s}_l) > k) : l = 1, \dots, 11760\}$. Lower left and lower right: SPOT function $T^{yr}(k)$ as a function of k linked to the map of $\{I(D^{yr}(\mathbf{s}_l) > k) : l = 1, \dots, 11760\}$. Units on the horizontal axis of both SPOT functions are degrees Celsius. The video is created by varying k from 1.0°C to 7.6°C in steps of 0.2°C .

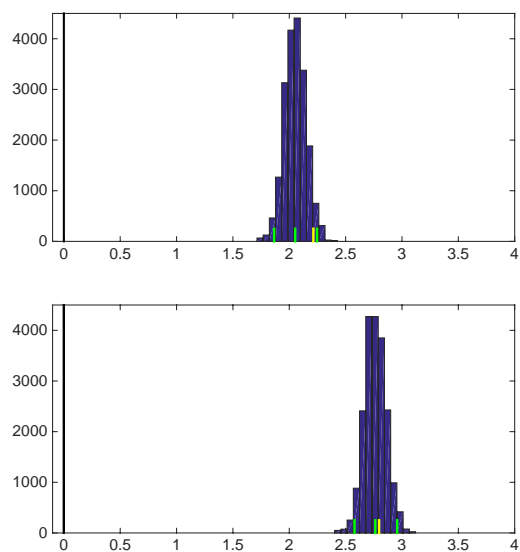


Fig. 3: The NARCCAP pixel \mathbf{s}_0 that contains NCAR's Mesa Lab is featured. Upper panel: Predictive distribution of $Y^{wi}(\mathbf{s}_0)$ given the data $\{D_{ij}(\cdot)\}$. Lower panel: Predictive distribution of $Y^{yr}(\mathbf{s}_0)$ given the data $\{D_{ij}(\cdot)\}$. For both plots, the vertical axis shows counts (out of 20,000) and units on the horizontal axis are degrees Celsius.

(Figure 4 video)

Fig. 4: Video showing inferential analysis of temperature change in degrees C in North America. Left panel: Map of PROT function $\{P^{wi}(k; \mathbf{s}_l) : l = 1, \dots, 11760\}$. Right panel: Map of PROT function $\{P^{yr}(k; \mathbf{s}_l) : l = 1, \dots, 11760\}$. Units on the horizontal axis of both PROT functions are degrees Celsius. The video is created by varying k from 1.0°C to 6.8°C in steps of 0.2°C .

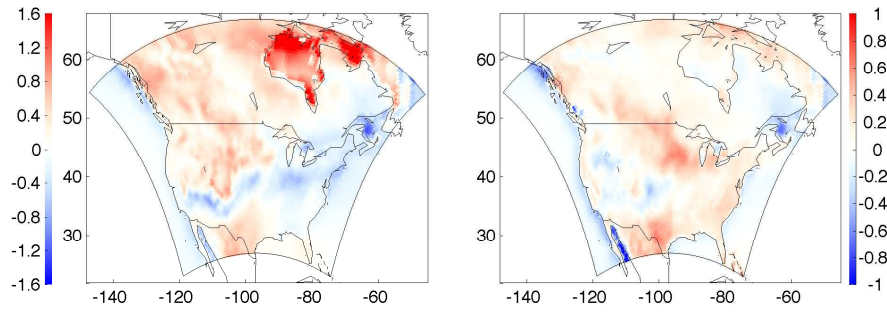


Fig. 5: Between-model variability of temperature change in degrees C in North America. Left panel: Map of $\{D_{14}(\mathbf{s}_l) - D_{24}(\mathbf{s}_l) : l = 1, \dots, 11760\}$ for the winter season. Right panel: Map of $\{(1/4) \sum_{j=1}^4 D_{1j}(\mathbf{s}_l) - (1/4) \sum_{j=1}^4 D_{2j}(\mathbf{s}_l) : l = 1, \dots, 11760\}$ for the whole year. The scale for both maps are in units of degrees Celsius.